



## Perceived prominence and scale types

Tøndering, John; Jensen, Christian

*Published in:*  
Proceedings Fonetik 2005

*Publication date:*  
2005

*Document version*  
Publisher's PDF, also known as Version of record

*Citation for published version (APA):*  
Tøndering, J., & Jensen, C. (2005). Perceived prominence and scale types. In A. Eriksson, & J. Lindh (Eds.), *Proceedings Fonetik 2005: the XVIIIth Swedish Phonetics Conference, May 25–27, 2005, Göteborg* (pp. 111-114)

## Perceived prominence and scale types

Christian Jensen<sup>1</sup> and John Tøndering<sup>2</sup>

<sup>1</sup>Department of English, Copenhagen Business School, Denmark

<sup>2</sup>Institute of Nordic Studies and Linguistics, Dept. of Linguistics, University of Copenhagen, Denmark

### Abstract

*Three different scales which have been used to measure perceived prominence are evaluated in a perceptual experiment. Average scores of raters using a multi-level (31-point) scale, a simple binary (2-point) scale and an intermediate 4-point scale are almost identical. The potentially finer gradation possible with the multi-level scale(s) is compensated for by having multiple listeners, which is also a requirement for obtaining reliable data. In other words, a high number of levels is neither a sufficient nor a necessary requirement. Overall the best results were obtained using the 4-point scale, and there seems to be little justification for using a 31-point scale.*

### Introduction

The purpose of this paper is to evaluate the use of different scales for measuring the perceived prominence of syllables and words. In this investigation only word-level prominence is considered.

Prominence, as perceived by groups of raters, has been measured on different types of scale: some use a 31-point scale from 0 to 30, first described in Fant & Kruckenberg (1989). The strength of this scale is that it allows for very fine gradation of the perceived prominence, even for a single rater, but this also makes the task quite difficult. Others, e.g. Wightman (1993), have proposed to use instead a simple binary (2-point) scale (0 or 1) and use the cumulative (or average) score of each word as an expression of its level of prominence, which results in much simpler task for the raters. The disadvantage of this simple scale is that it may force raters to conflate items which they perceive as “different, but within the same category”, which could lead to a reduced or lost ability to distinguish variations in perceived prominence at either end of the prominence continuum. For example accented words with or without special emphasis. In addition, the level of gradation you achieve with this scale is directly proportional to the number of raters: to get the same gradation as is (potentially) possible with the scale from 0 to 31 you need 30 raters. As a possible compromise between these two scales one could use a 4-point scale (e.g. from 0 to 3). While this scale is much simpler

than the 31-point scale it still allows raters to make some gradation in their prominence evaluations.

We investigated the three prominence scales outlined above with the purpose of answering two overall questions: does the choice of scale influence the results with regard to 1) the perceived prominence relations of words in utterances, and 2) the ability to make observations about statistically significant differences between words. These questions were addressed from the point of view of three relevant linguistic parameters which are known to be associated with perceived prominence: *part of speech membership*, *information structure* and *correlation with  $F_0$* .

### Method

The speech material chosen to evaluate the scales was two short monologues from the Danish DanPASS project (<http://www.cphling.dk/pers/ng/danpass.htm>), both recordings of a map task activity. The two monologues, by two different male speakers, included a total of 123 words. The monologues were divided into shorter phrases which were presented via a web page (one phrase per page). The raters could hear the phrase as many times as they wanted by pressing a “play” button, and indicated their judgment by clicking the appropriate scale point. Time consumption and a count of sound file playbacks were recorded for each phrase.

A large group of raters participated in the experiment and were randomly assigned to a specific scale. Equally sized groups of 19 raters (the size of the smallest group) were selected for the analyses. The instructions to the raters were presented from the web page and were identical for all three groups, except for the details about the specific scale. The concept of prominence was explained and exemplified, and raters were advised that prominence might be a question of “more or less”. 0 represented *no prominence*, but no other scale points were defined. Prominent words could be assigned values *up to* the scale maximum. Raters using the 2-point scale were informed that they could not grade their ratings but were given a forced choice.

## Results

### Reliability

Note: the phrase “the 2/4/31-point scale” is used in the following as shorthand expressions of “the prominence ratings obtained from the group of listeners using the 2/4/31-point scale”.

The reliability of the data was tested by calculating Cronbach’s  $\alpha$  coefficient, which expresses the extent to which the scores of the individual raters covary. The coefficients for all three groups are high (from 0.94 to 0.96) and the difference between them is nonsignificant ( $M = 1.02$ ,  $p > 0.05$ ).

### Comparison of prominence ratings

The first question to be addressed is whether the prominence ratings on the three scales express the same relations between words. In order to be able to make direct comparisons all scores were normalised by dividing each value with the scale maximum (1, 3 or 30, respectively), which fits all data to a normalised scale of 0 to 1 without affecting the relations between scores. These values were then plotted on a line chart for simple visual inspection. An example diagram of one phrase is shown in Fig. 1.

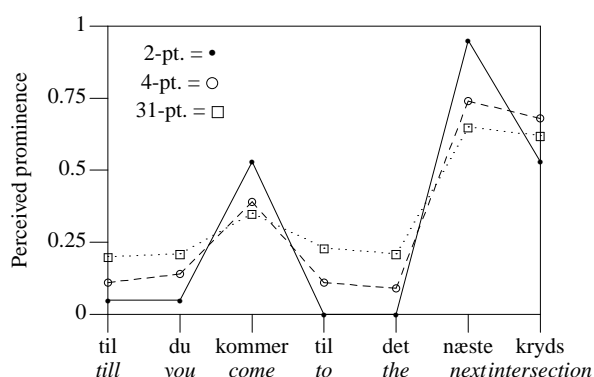


Figure 1: Prominence of selected phrase – all scales

The diagrams showed a high level of agreement across the three scales, which was further tested in a correlation analysis (Spearman’s  $\rho$ ). The result can be seen in Table 1.

Table 1: Correlation coefficients (Spearman’s  $\rho$ ) across all three scales

Correlation	4-pt	31-pt
2-pt	0.933	0.926
4-pt	—	0.964

The correlation coefficients were high for each scale pair and quite similar, with the best correlation apparently between the 4-point scale

and the 31-point scale. The preliminary conclusion is clear: raters arrive at approximately the same rank order of perceived prominence regardless of the scale used.

It appears from Fig. 1 that the 2-point scale displays somewhat larger variation in values between the scale minimum and maximum than the 4-point scale and especially the 31-point scale. This was in fact a general trend demonstrating a certain compression of values on the 31-point scale (and to a lesser degree the 4-point scale), while the 2-point scale has more mean values near the scale extremes. Analyses of the distribution of scores (inter-quartile range for each rater and visual inspection of  $x$ - $y$  plots) showed that many raters on the 31-point scale assigned most ratings to a restricted – sometimes very restricted – range of the scale, either at the lower, the middle or the higher end of the scale. There are therefore no *mean values* at the scale extremes, although there were many individual scores near the minimum and maximum values.

### Obtaining significant differences

One very important aspect of choosing a scale is whether it will affect the ability to obtain statistically significant differences between test items. The hypothesis might be that scales with too few points (most notably the 2-point scale) would mask subtle perceptual differences which could be brought out with more scale points.

This suitability of the three scales for quantitative analysis was tested by examining the association between perceived prominence and three linguistic phenomena: part of speech membership, information structure and a specific acoustic correlate, namely  $F_0$ . The purpose was to see if the data obtained by using three different scales will lead to different conclusions about linguistic structure.

### Comment on the statistical procedures

Since it is not possible to compare results directly across scale types we simply decided to use the statistical procedures which were felt to be most appropriate for each individual scale. This resembles quite well the choice which researchers would be forced to make when they are making a choice about scale type.

For all scales we have decided to use non-parametric methods. For significance testing on the 2-point scale we use the Fisher exact test or a chi-square test with corrections for continuity (when  $n > 40$ ), and for the other two scales we use the Wilcoxon-Mann-Whitney test with correction for ties (WMW).

Table 2: Prominence ratings and parts of speech. Left braces indicate non-significant differences. Non-adjacent, nonsignificant differences on the 31-pt scale: *adv-v*, *art-prep*

<i>Scale →</i>		<i>n</i>	2-point		4-point		31-point	
<i>Part of speech</i>			<i>Ranked</i>	$\bar{x}$	<i>Ranked</i>	$\bar{x}$	<i>Ranked</i>	$\bar{x}$
1	Adjectives	9	<i>adj</i>	0.92	<i>adj</i>	0.73	<i>adj</i>	0.67
2	Nouns	28	<i>n</i>	0.78	<i>n</i>	0.66	<i>n</i>	{0.63
3	Interjections	3	<i>int</i>	{0.60	<i>int</i>	0.50	<i>int</i>	{0.58
4	Adverbs	12	<i>adv</i>	{0.58	<i>adv</i>	0.38	<i>adv</i>	0.40
5	Verbs	13	<i>v</i>	{0.34	<i>v</i>	{0.30	<i>pron</i>	{0.35
6	Pronouns	16	<i>pron</i>	{0.33	<i>pron</i>	{0.30	<i>v</i>	{0.35
7	Conjunctions	10	<i>conj</i>	{0.17	<i>prep</i>	0.21	<i>prep</i>	0.28
8	Articles	2	<i>art</i>	{0.13	<i>conj</i>	{0.13	<i>conj</i>	{0.24
9	Prepositions	30	<i>prep</i>	{0.10	<i>art</i>	{0.12	<i>art</i>	{0.22

### Parts of speech

The mean prominence ratings of nine parts of speech are listed in Table 2, ordered according to their ranking on each scale. These ranking are very similar for all three scales. The only difference which can be detected is the relegation of *prepositions* to ninth place on the 2-point scale, instead of the seventh place it holds on the other two scales. (The different ranking of *pronouns* and *verbs* on the 31-point scale is irrelevant.) Most of the differences between the classes are significant: except for two cases on the 31-point scale (see the table caption) all differences between classes which are not adjacent in the rankings are significant, and of the differences between adjacent classes four are nonsignificant on the 2-point scale, two are nonsignificant on the 4-point scale, and three are nonsignificant on the 31-point scale (giving a total of five differences which are not significant for this scale). These figures are quite similar, with a small bias in favour of the 4-point scale, where the highest number of significant differences was found.

### Information structure

Chafe (1994) states that new information is more prominent than non-new information. To test the validity of this statement we compared the prominence ratings of all words carrying new information with the most prominent word carrying non-new information in the same phrase (20 cases), thus testing the hypothesis that new information is more prominent than other information ( $H_1$ ).  $H_0$  states that the perceived prominence of the new information is less than or equal to that of the given/accessible information.

In four cases (three on the 31-point scale) the new information is not more prominent than the non-new information, in which case  $H_0$  cannot

be dismissed. Of the remaining 16 (17) cases, where the new information had higher prominence ratings than the non-new information, nine were significant on the 2-point scale (Fisher exact test, one-tailed,  $p < 0.05$ ); 15 were significant on the 4-point scale and 14 on the 31-point scale (WMW, one-tailed,  $p < 0.05$ ).

Here we find a clear difference between the 2-point scale and the 4-point and 31-point scales in the number of significant differences. Our conclusion about the relative prominence levels of new versus non-new information would therefore be affected by our choice of scale, provided that we want to verify observed differences in mean ratings statistically.

### Correlation with $F_0$

The prominence level of a Danish accented syllable, and of the word in which it occurs, is generally felt to be associated with, among other cues, a rise in  $F_0$ . The greater the rise, the more prominent the syllable is perceived to be. For this investigation two  $F_0$  values were measured for all words in which such a rise occurs: the  $F_0$  trough and the  $F_0$  peak value within the domain of onset of the accented vowel and the end of the word (since we were concerned with word level prominence). The rise is expressed as the difference in semitones between these two values, and the values for the rises were then correlated against the prominence ratings from the three scales. The results are displayed in Table 3.

The correlation coefficients are very similar for the three data sets, indicating the the association between prominence and  $F_0$  can be described equally well regardless of the scale used. To the (slight) extent that any difference can be detected it seems that the correlation is better with data obtained on the 4-point scale.

Table 3: Correlation (Spearman's  $\rho$ ) between perceived prominence and  $F_0$

Scale	$\rho$
2-pt	0.593
4-pt	0.626
31-pt	0.606

### Rater effort, or level of difficulty

In a few places we have described the 2-point scale, and to some extent the 4-point scale, as “simpler” and less difficult for the rater than the 31-point scale. At least this was our expectation, and as an attempt to capture this we measured the time consumption for each phrase and number of times the raters listened to each phrase. The hypothesis is that both of these measures will increase with an increase in the number of scale points.

This hypothesis was in fact borne out: there is an increase in time consumption of 18% when going from two to four scale points, and an increase of 42% when going from two to 31 points. All pairwise comparisons between the three scales are significant (t-tests, one-tailed,  $p < 0.05$ ). The pattern is less clear for the number of playbacks, where only the tendency for more playbacks on the 31-point scale compared with the 2- and 4-point scales is statistically significant.

It must be concluded, though, that using more scale points will result in a somewhat higher “cost”.

### Discussion and conclusion

Two main questions were asked about the influence of scale type on ratings of perceived prominence: 1) do we get the same prominence relations in utterances, as expressed in mean values and rankings, and 2) does scale type affect our ability to make observations about statistically significant differences between words. The overall conclusion must be that the perceived prominence relations in the utterances are very similar whether expressed on a 2-point scale, a 4-point scale or a 31-point scale. The differences are small and are mostly caused by a tendency for some raters to prefer a restricted range within a multi-level scale. The differences are also relatively small when it comes to statistical testing of observations, but it does seem that raising the number of scale points from two to four yields slightly better results: there are more significant differences between the part of

speech categories and between words with new versus given/accessible information, and the correlation with  $F_0$  is better. No such improvement can be obtained, however, by raising the number of scale point to 31. On the contrary we find slightly fewer significant differences on this scale.

One reason for this finding may be that it is too difficult for untrained listeners to use the 31-point scale. In a parallel experiment (to be reported elsewhere) we had five expert listeners rate the same phrases as in this experiment (with slightly different instructions). The performance of this group was generally better than any random group of five untrained listeners (higher Cronbach  $\alpha$  coefficient and more significant differences), which indicates that they did in fact do better on this scale. The analysis also showed, however, that five expert listeners cannot replace a larger group of untrained listeners if the objective is to find statistically significant differences – the number of observations becomes too small.

It was shown that “expenses”, in terms of especially time consumption, grew with an increase in the number of scale points. Combined with the above observations this points to a recommendation of using many listeners rating on a scale with relatively few levels. A 2-point scale may then be adequate for most purposes and makes for the simplest and fastest task, but it would appear that increasing the number of levels to four results in slightly better performance. There seems to be no justification for using a 31-point scale, unless the requirement of using many listeners cannot be met. The task becomes more difficult and takes more time, and there is no gain in terms of precision or “discriminatory power” to balance the extra cost.

### References

- Fant, G. and Kruckenberg, A., “Preliminaries to the study of Swedish prose reading and reading style”, STL-QPSR 2/1989:1–83, 1989.
- Wightman, C., “Perception of multiple levels of prominence in spontaneous speech”, ASA 126th Meeting Denver 1993 (abstract).
- Chafe, W., “Discourse consciousness, and time: the flow and displacement of conscious experience in speaking and writing”, The University of Chicago Press, 1994.